

Optimization

Lectured by Dr Dante Kalise

Typed by Aris Zhu Yi Qing

April 28, 2022

Contents

1	Mathematical Preliminaries	3
1.1	Topological Concepts	3
1.2	Multi-variable Calculus	4
1.3	Positive Definiteness of Matrix	4
2	Unconstrained Optimization	6
2.1	Optimums	6
2.2	Second-order Optimality Conditions	6
2.3	Attainment of Minimal/Maximal Points	7
2.4	Global Optimality Conditions	7
3	Linear Least Squares	8
3.1	Problem Formulation	8
3.2	Data Fitting	8
3.3	Regularized Least Squares	9
3.4	Denoising	9
4	The Gradient Method	10
4.1	Descent Direction	10
4.2	Stepsize Selection Rules	11
4.3	Convergence	11
4.4	Condition Number and Convergence for Quadratic Function	12
4.5	Scaled Gradient Method	13
4.6	The Kaczmarz Algorithm	14
4.7	Stochastic Gradient Descent	14
5	Convexity	16
5.1	Convex Sets	16
5.2	Convex Hull	17
5.3	Convex Functions	19
5.4	First-order Characterization of Convex Functions	20
5.5	Second-order Characterization of Convex Functions	20
5.6	More Results of Convex Function	21

6	Convex Optimization	23
6.1	Problem Definition	23
6.2	Stationarity	24
6.3	Orthogonal Projection Operator	25
6.4	Gradient Projection Method	26
6.5	Separation Theorem	26
6.6	KKT Conditions	27
6.7	Orthogonal projections	28
7	Duality	31
7.1	The Primal and Dual Problems	31
7.2	Weak and Strong Duality	32

Chapter 1

Mathematical Preliminaries

1.1 Topological Concepts

Definition 1. The **open ball** with center $c \in \mathbb{R}^n$ and radius r is

$$B(c, r) = \{\mathbf{x} : \|\mathbf{x} - c\| < r\}.$$

Similarly, the **closed ball** with center c and radius r is

$$B[c, r] = \{\mathbf{x} : \|\mathbf{x} - c\| \leq r\}.$$

Definition 2. Given a set $U \subseteq \mathbb{R}^n$, a point $\mathbf{c} \in U$ is called an **interior point** of U if $\exists r > 0$ for which $B(\mathbf{c}, r) \subseteq U$. The set of all interior points of a given set U is called the interior of the set and is denoted by

$$\text{int}(U) = \{\mathbf{x} \in U : B(\mathbf{x}, r) \subseteq U \text{ for some } r > 0\}.$$

Definition 3. Given a set $U \subseteq \mathbb{R}^n$, a **boundary point** of U is a vector $\mathbf{x} \in \mathbb{R}^n$ satisfying that any neighbourhood of \mathbf{x} contains at least one point in U and at least one point in its complement U^c . We denote

$$\text{bd}(U) = \text{The set of all boundary points of a set } U.$$

Definition 4. The **closure** of a set $U \subseteq \mathbb{R}^n$ is the smallest closed set containing U , denoted by $\text{cl}(U)$ with

$$\text{cl}(U) = U \cup \text{bd}(U).$$

Definition 5. A set $U \subseteq \mathbb{R}^n$ is called **bounded** if $\exists M > 0$ for which $U \subseteq B(0, M)$.

Definition 6. A set $U \subseteq \mathbb{R}^n$ is called **compact** if it is closed and bounded.

1.2 Multi-variable Calculus

Definition 7. The **directional derivative** of a scalar function f w.r.t. \mathbf{d} at a point \mathbf{x} is denoted as

$$f'(\mathbf{x}; \mathbf{d}) = \nabla f(\mathbf{x})^T \mathbf{d}$$

Theorem 8. Given the general quadratic functions of the form

$$f(\mathbf{w}) = \mathbf{w}^T A \mathbf{w} + \mathbf{b}^T \mathbf{w} + \gamma$$

we have

$$\nabla f(\mathbf{w}) = (A^T + A)\mathbf{w} + \mathbf{b}, \quad \nabla^2 f(\mathbf{w}) = A + A^T.$$

If A is symmetric, then

$$\nabla f(\mathbf{w}) = 2A\mathbf{w} + \mathbf{b}, \quad \nabla^2 f(\mathbf{w}) = 2A.$$

1.3 Positive Definiteness of Matrix

Proposition 9. Let A be a positive definite (semidefinite) matrix, then

- the diagonal elements of A are positive (nonnegative)
- $\text{Tr}(A)$ and $\det(A)$ are positive (nonnegative)

(Test 1) Theorem 10. Let $A \in \mathbb{R}^{n \times n}$ be symmetric, then

- A is positive definite (semidefinite) iff all its eigenvalues are positive (nonnegative).
- A is indefinite iff it has at least one positive eigenvalue and at least one negative eigenvalue.

Definition 11. Let $A \in \mathbb{R}^{n \times n}$ be symmetric, then

- A is **diagonally dominant** if

$$|A_{ii}| \geq \sum_{j \neq i} |A_{ij}| \quad \forall i = 1, 2, \dots, n$$

- A is **strictly diagonally dominant** if

$$|A_{ii}| > \sum_{j \neq i} |A_{ij}| \quad \forall i = 1, 2, \dots, n$$

(Test 2) Theorem 12. If $A \in \mathbb{R}^{n \times n}$ is symmetric, diagonally dominant with positive (nonnegative) diagonal elements, then A is positive definite (semidefinite).

Chapter 2

Unconstrained Optimization

2.1 Optimums

Definition 13. Let $f : S \rightarrow \mathbb{R}$ be defined on a set $S \subseteq \mathbb{R}^n$, then $\forall \mathbf{x} \in S$,

$\mathbf{x}^* \in S$ is a **global minimum** point of f over S if $f(\mathbf{x}) \geq f(\mathbf{x}^*)$,

$\mathbf{x}^* \in S$ is a **strict global minimum** point of f over S if $f(\mathbf{x}) > f(\mathbf{x}^*)$,

and similar definitions for maximum.

Definition 14. Let $f : S \rightarrow \mathbb{R}$ be defined on a set $S \subseteq \mathbb{R}^n$, $\mathbf{x}^* \in S$ is a **local minimum** of f over S if $\exists r > 0$ s.t. $f(\mathbf{x}^*) \leq f(\mathbf{x})$ for any $\mathbf{x} \in S \cap B(\mathbf{x}^*, r)$. Similar definitions for **strict local minimum** and maximum.

Definition 15. Let $f : U \rightarrow \mathbb{R}$ be a function defined on a set $U \subseteq \mathbb{R}^n$. Suppose that $\mathbf{x}^* \in \text{int}(U)$ and that all the partial derivatives of f are defined at \mathbf{x}^* , then \mathbf{x}^* is called a **stationary point** of f if $\nabla f(\mathbf{x}^*) = 0$.

2.2 Second-order Optimality Conditions

Theorem 16. Let $f : U \rightarrow \mathbb{R}$ be a function defined on an open set $U \subseteq \mathbb{R}^n$. Suppose that f is twice continuously differentiable over U and that \mathbf{x}^* is a stationary point, then

- \mathbf{x}^* is a local minimum point $\iff \nabla^2 f(\mathbf{x}^*) \succeq 0$.
- \mathbf{x}^* is a strict local minimum point $\iff \nabla^2 f(\mathbf{x}^*) \succ 0$.
- similar necessary and sufficient conditions for (strict) local maximum point

Definition 17. Let $f : U \rightarrow \mathbb{R}$ be a continuously differentiable function defined on an open set $U \subseteq \mathbb{R}^n$. A stationary point $\mathbf{x}^* \in U$ is called a **saddle point** of f over U if it is neither a local minimum nor a local maximum point of f over U .

Theorem 18. Let $f : U \rightarrow \mathbb{R}$ be a continuously differentiable function defined on an open set $U \subseteq \mathbb{R}^n$. Suppose that f is twice continuously differentiable over U and that \mathbf{x}^* is a stationary point. Then

$$\nabla^2 f(\mathbf{x}^*) \text{ is an indefinite matrix} \implies \mathbf{x}^* \text{ is a saddle point of } f \text{ over } U.$$

2.3 Attainment of Minimal/Maximal Points

(Weierstrass') Theorem 19. Let f be a continuous function defined over a nonempty compact set $C \subseteq \mathbb{R}^n$. Then \exists a global minimum point of f over C and a global maximum point of f over C .

Definition 20. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous function over \mathbb{R}^n . f is called **coercive** if

$$\lim_{\|\mathbf{x}\| \rightarrow \infty} f(\mathbf{x}) = \infty$$

Theorem 21. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous and coercive function and let $S \subseteq \mathbb{R}^n$ be a nonempty closed set. Then f attains a global minimum point on S .

2.4 Global Optimality Conditions

Theorem 22. Let f be a twice continuously differentiable function defined over \mathbb{R}^n . Let $\mathbf{x}^* \in \mathbb{R}^n$ be a stationary point of f . Then

$$\nabla^2 f(\mathbf{x}) \succeq 0 \quad \forall \mathbf{x} \in \mathbb{R}^n \implies \mathbf{x}^* \text{ is a global minimum point of } f.$$

Proposition 23. Let $f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c$, with $A \in \mathbb{R}^{n \times n}$ symmetric, then

1. \mathbf{x} is a stationary point of f iff $A\mathbf{x} = -\mathbf{b}$.
2. if $A \succeq 0$, then \mathbf{x} is a global minimum point of f iff $A\mathbf{x} = -\mathbf{b}$.
3. if $A \succ 0$, then $\mathbf{x} = -A^{-1}\mathbf{b}$ is a strict global minimum point of f .

Chapter 3

Linear Least Squares

3.1 Problem Formulation

Consider the linear system

$$S\mathbf{x} \approx \mathbf{b}, \quad (S \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m, m > n)$$

To solve the above system, the usual approach is to transform it to become

$$\min_{\mathbf{x}} \|S\mathbf{x} - \mathbf{b}\|^2 \iff \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ f(\mathbf{x}) \equiv \mathbf{x}^T S^T S \mathbf{x} - 2\mathbf{b}^T S \mathbf{x} + \|\mathbf{b}\|^2 \right\}.$$

Note that $\nabla^2 f(\mathbf{x}) = 2S^T S \succeq 0$ since $\mathbf{x}^T S^T S \mathbf{x} = (S\mathbf{x})^T (S\mathbf{x}) = \|S\mathbf{x}\|^2 \geq 0$. Therefore, the unique optimal solution \mathbf{x}_{LS} is the solution $\nabla f(\mathbf{x}) = 0$, namely

$$(S^T S)\mathbf{x}_{\text{LS}} = S^T \mathbf{b} \implies \mathbf{x}_{\text{LS}} = (S^T S)^{-1} S^T \mathbf{b}.$$

3.2 Data Fitting

1. For dataset (\mathbf{s}_i, b_i) where $\mathbf{s}_i \in \mathbb{R}^n$ and $b_i \in \mathbb{R}$, we could transform to problem

$$\min_{\mathbf{x}} \sum_{i=1}^m (\mathbf{s}_i^T \mathbf{x} - b_i)^2 \implies \min_{\mathbf{x}} \|S\mathbf{x} - \mathbf{b}\|^2$$

2. For polynomial fitting, given a set of points $\mathbb{R}^2 : (u_i, y_i)$, the associated linear system is

$$\begin{pmatrix} 1 & u_1 & u_1^2 & \cdots & u_1^d \\ 1 & u_2 & u_2^2 & \cdots & u_2^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & u_m & u_m^2 & \cdots & u_m^d \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_m \end{pmatrix}$$

3.3 Regularized Least Squares

A Regularized Least Square problem is formulated as

$$\min_{\mathbf{x}} \|S\mathbf{x} - \mathbf{b}\|^2 + \lambda R(\mathbf{x}),$$

where λ is the regularization parameter and $R(\cdot)$ is the regularization function (also called a *penalty* function). A common choice is a quadratic regularization function:

$$\min_{\mathbf{x}} \|S\mathbf{x} - \mathbf{b}\|^2 + \lambda \|D\mathbf{x}\|^2$$

with its optimal solution being

$$\mathbf{x}_{\text{RLS}} = (S^T S + \lambda D^T D)^{-1} S^T \mathbf{b}$$

since $\nabla f = 2S^T S\mathbf{x} - 2S^T \mathbf{b} + 2\lambda D^T D\mathbf{x} = 0$.

3.4 Denoising

Suppose a noisy measurement of a signal $\mathbf{x} \in \mathbb{R}^n$ is given

$$\mathbf{b} = \mathbf{x} + \mathbf{w}$$

where \mathbf{x} is the “true” unknown signal, \mathbf{w} is the unknown noise and \mathbf{b} is the (known) measures vector. We could define

$$R(\mathbf{x}) = \|L\mathbf{x}\|^2, \text{ where } L = \begin{pmatrix} 1 & -1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 1 & -1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & -1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 & -1 \end{pmatrix}$$

as the regularization function to penalize any sudden variations in signal. The RLS is thus

$$\min_{\mathbf{x}} \|\mathbf{x} - \mathbf{b}\|^2 + \lambda \|L\mathbf{x}\|^2$$

with its direct solution being

$$\mathbf{x}_{\text{RLS}}(\lambda) = (I + \lambda L^T L)^{-1} \mathbf{b}.$$

Chapter 4

The Gradient Method

4.1 Descent Direction

Definition 24. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function. A vector $\mathbf{0} \neq \mathbf{d} \in \mathbb{R}^n$ is called a **descent direction** of f at \mathbf{x} if

$$f'(\mathbf{x}; \mathbf{d}) = \nabla f(\mathbf{x})^T \mathbf{d} < 0.$$

Example 25. The descent direction can be $\mathbf{d} = -\nabla f(\mathbf{x})$, since as long as $\nabla f(\mathbf{x}) \neq \mathbf{0}$ (\mathbf{x} is a non-stationary point), we have

$$f'(\mathbf{x}; -\nabla f(\mathbf{x})) = -\nabla f(\mathbf{x})^T \nabla f(\mathbf{x}) = -\|\nabla f(\mathbf{x})\|^2 < 0.$$

Lemma 26. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function. Let $\mathbf{x} \in \mathbb{R}^n$. Suppose that \mathbf{d} is a descent direction of f at \mathbf{x} , then

$$\exists \epsilon > 0 \text{ s.t. } \forall t \in (0, \epsilon], f(\mathbf{x} + t\mathbf{d}) < f(\mathbf{x}).$$

Lemma 27. Let f be a continuously differentiable function and $\mathbf{x} \in \mathbb{R}^n$ be a non-stationary point ($\nabla f(\mathbf{x}) \neq \mathbf{0}$), then the optimal solution of

$$\min_{\mathbf{d}} \{f'(\mathbf{x}; \mathbf{d}) : \|\mathbf{d}\| = 1\}$$

is $\mathbf{d} = -\frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|}$.

Lemma 28. Let $\{\mathbf{x}^k\}_{k \geq 0}$ be the sequence generated by the gradient descent method with *exact* line search for solving a problem of minimizing a continuously differentiable function f . Then $\forall k = 0, 1, 2, \dots$,

$$(\mathbf{x}^{k+2} - \mathbf{x}^{k+1})^T (\mathbf{x}^{k+1} - \mathbf{x}^k) = 0.$$

4.2 Stepsize Selection Rules

Finding the right $t^k \in \mathbb{R}^n$, called the **stepsize**, is referred in the literature as **line search**.

1. Constant stepsize: $t^k = \bar{t} \forall k$.
2. Exact stepsize: t^k is a minimizer of f along the ray $\mathbf{x}_t^k \mathbf{d}^k$:

$$t^k \in \operatorname{argmin}_{t \geq 0} f(\mathbf{x}^k + t \mathbf{d}^k)$$

3. Backtracking (Armijo rule): let $s > 0, \alpha \in (0, 1), \beta \in (0, 1)$, and initial stepsize $t^k = s$, while

$$f(\mathbf{x}^k) - f(\mathbf{x}^k + t^k \mathbf{d}^k) < -\alpha t^k \nabla f(\mathbf{x}^k)^T \mathbf{d}^k$$

set $t^k := \beta t^k$, iterating until achieving the **sufficient decrease property**

$$f(\mathbf{x}^k) - f(\mathbf{x}^k + t^k \mathbf{d}^k) \geq -\alpha t^k \nabla f(\mathbf{x}^k)^T \mathbf{d}^k.$$

4.3 Convergence

Definition 29. Let f be a continuously differentiable function over \mathbb{R}^n . We say that f has a **Lipschitz gradient** if

$$\exists L \geq 0 \text{ s.t. } \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|.$$

L is called the **Lipschitz constant**.

Comments:

- The class of functions with Lipschitz gradient with constant L is denoted as $C_L^{1,1}(\mathbb{R}^n)$ or just $C_L^{1,1}$. When L is irrelevant, we simply denote the class by $C^{1,1}$.
- If ∇f is Lipschitz with constant L , then it is also Lipschitz with constant $L' \forall L' \geq L$.
- Linear functions: Given $\mathbf{a} \in \mathbb{R}^n$, the function $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$ is in $C_0^{1,1}$.
- Quadratic functions: Let $A \in \mathbb{R}^{n \times n}$, $\mathbf{b} \in \mathbb{R}^n$, and $c \in \mathbb{R}$, then the function $f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c$ is $C_{2\|A\|_2}^{1,1}$.

Theorem 30. Let f be a continuously differentiable function over \mathbb{R}^n . Then

$$f \in C_L^{1,1}(\mathbb{R}^n) \iff \|\nabla^2 f(\mathbf{x})\| \leq L \forall \mathbf{x} \in \mathbb{R}^n.$$

(Sufficient decrease of the gradient method) Lemma 31. Let $f \in C_L^{1,1}(\mathbb{R}^n)$. Let $\{\mathbf{x}^k\}_{k \geq 0}$ be the sequence generated by the gradient method for solving

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

with one of the following stepsize strategies:

- constant stepsize $\bar{t} \in (0, \frac{2}{L})$,
- exact line search
- backtracking procedure with parameters $s \in \mathbb{R}_{++}$, $\alpha \in (0, 1)$, and $\beta \in (0, 1)$,

then

$$f(\mathbf{x}^k) - f(\mathbf{x}^{k+1}) \geq M \left\| \nabla f(\mathbf{x}^k) \right\|^2$$

where

$$M = \begin{cases} \bar{t} \left(1 - \frac{\bar{t}L}{2}\right) & \text{constant stepsize} \\ \frac{1}{2L} & \text{exact line search} \\ \alpha \min \left\{ s, \frac{2(1-\alpha)\beta}{L} \right\} & \text{backtracking} \end{cases}$$

(Convergence of the gradient method) Theorem 32. Let $f \in C_L^{1,1}(\mathbb{R}^n)$ and is bounded below over \mathbb{R}^n . Let $\{\mathbf{x}^k\}_{k \geq 0}$ be the sequence generated by the gradient method for solving

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

with one of the following stepsize strategies:

- constant stepsize $\bar{t} \in (0, \frac{2}{L})$,
- exact line search
- backtracking procedure with parameters $s \in \mathbb{R}_{++}$, $\alpha \in (0, 1)$, and $\beta \in (0, 1)$,

then

1. $\forall k, f(\mathbf{x}^{k+1}) < f(\mathbf{x}^k)$ unless $\nabla f(\mathbf{x}^k) = 0$.
2. $\nabla f(\mathbf{x}^k) \rightarrow 0$ as $k \rightarrow \infty$.

4.4 Condition Number and Convergence for Quadratic Function

Definition 33. Let $A \in \mathbb{R}^{n \times n}$ be positive definite, Then the **condition number** of A is

$$\kappa(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$$

where $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ are the largest and smallest eigenvalues respectively.

(Kantorovich inequality) Lemma 34. Let $A \in \mathbb{R}^{n \times n}$ be positive definite. Then

$$\forall \mathbf{0} \neq \mathbf{x} \in \mathbb{R}^n, \frac{(\mathbf{x}^T \mathbf{x})^2}{(\mathbf{x}^T A \mathbf{x})(\mathbf{x}^T A^{-1} \mathbf{x})} \geq \frac{4\lambda_{\max}(A)\lambda_{\min}(A)}{(\lambda_{\max}(A) + \lambda_{\min}(A))^2}.$$

(Convergence for quadratic function) Theorem 35. Let $\{\mathbf{x}^k\}_{k \geq 0}$ be the sequence generated by the gradient method for solving

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}^T A \mathbf{x}) \quad (A \succ 0),$$

then $\forall k = 0, 1, \dots,$

$$f(\mathbf{x}^{k+1}) \leq \left(\frac{\lambda_{\max}(A) - \lambda_{\min}(A)}{\lambda_{\max}(A) + \lambda_{\min}(A)} \right)^2 f(\mathbf{x}^k) = \left(\frac{\kappa(A) - 1}{\kappa(A) + 1} \right)^2 f(\mathbf{x}^k).$$

4.5 Scaled Gradient Method

A way to mitigate the slow convergence due to poor conditioning of the Hessian is to formulate a rescaled version of the problem. From the minimization problem

$$\min \{f(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\}$$

we introduce a nonsingular matrix $S \in \mathbb{R}^{n \times n}$ to make the linear change of variables $\mathbf{x} = S\mathbf{y}$ and obtain the equivalent problem

$$\min \{g(\mathbf{y}) \equiv f(S\mathbf{y}) : \mathbf{y} \in \mathbb{R}^n\}$$

Since $\nabla g(\mathbf{y}) = S^T \nabla f(S\mathbf{y}) = S^T \nabla f(\mathbf{x})$, the gradient method for the rescaled problem reads

$$\mathbf{y}^{k+1} = \mathbf{y}^k - t^k S^T \nabla f(S\mathbf{y}^k).$$

Multiplying both sides by S , with $\mathbf{x}^k = S\mathbf{y}^k$, and define $D = SS^T$, we have

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t^k D \nabla f(\mathbf{x}^k).$$

Since $D \succ 0$, so

$$f'(\mathbf{x}^k; -D \nabla f(\mathbf{x}^k)) = -\nabla f(\mathbf{x}^k)^T D \nabla f(\mathbf{x}^k) < 0.$$

A well-known choice for D^k is to pick $D^k = (\nabla^2 f(\mathbf{x}^k))^{-1}$ (Newton's method). Another alternative is to use a diagonal scaling, e.g.

$$(D^k)_{ii} = \left(\frac{\partial^2 f(\mathbf{x}^k)}{\partial x_i^2} \right)^{-1}$$

4.6 The Kaczmarz Algorithm

The *Kaczmarz Algorithm* solves the linear system

$$A\mathbf{x} = \mathbf{b}$$

by iterating projections along the i -th row of the matrix A , denoted by \mathbf{a}_i^T :

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \frac{b_i - \mathbf{a}_i^T \mathbf{x}^k}{\|\mathbf{a}_i\|^2} \mathbf{a}_i$$

In the original Kaczmarz algorithm, the i -th row is chosen periodically by cycling through all rows. If chooses i -th row randomly, we can show that the algorithm converges exponentially, and this is known as *randomized Kaczmarz Algorithm*.

The algorithm works because the problem of solving the linear system $A\mathbf{x} = \mathbf{b}$ could be formulated as an optimization problem

$$\min_{\mathbf{x}} \frac{1}{2m} \|A\mathbf{x} - \mathbf{b}\|^2 = \frac{1}{2m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{x} - b_i)^2$$

for which the gradient descent method could be constructed as

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{t}{m} A^T (A\mathbf{x} - \mathbf{b})$$

but the problem could also be formulated as

$$\min_{\mathbf{x}} \frac{1}{2m} \|A\mathbf{x} - \mathbf{b}\|^2 = \frac{1}{2m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{x} - b_i)^2 = \frac{1}{2} \mathbb{E}_i [\mathbf{a}_i^T \mathbf{x} - b_i]^2,$$

which can then be translated to the action of randomly picking a row of A , becoming

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{t}{m} (\mathbf{a}_i^T \mathbf{x} - b_i) \mathbf{a}_i$$

4.7 Stochastic Gradient Descent

Theorem 36. Assuming that

- The cost $g(\mathbf{x})$ is such that

$$\|\nabla g(\mathbf{x}) - \nabla g(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \text{and} \quad \nabla^2 g(\mathbf{x}) \succeq \mu I.$$

- The sample gradient $\nabla Q_i(\mathbf{x}^k)$ is an unbiased estimate of $\nabla g(\mathbf{x}^k)$.

-

$$\forall \mathbf{x}, \mathbb{E}_i \left[\|\nabla Q_i(\mathbf{x})\|^2 \right] \leq \sigma^2 + c\|\nabla g(\mathbf{x})\|^2.$$

Then if $t^k \equiv t \leq \frac{1}{Lc}$, then SGD achieves

$$\mathbb{E} \left[g(\mathbf{x}^k) - g(\mathbf{x}^*) \right] \leq \frac{tL\sigma^2}{2\mu} + (1 - t\mu)^k (g(\mathbf{x}^0) - g(\mathbf{x}^*)).$$

Comments

1. Fast (linear) convergence during the first iterations.
2. Convergence to a neighbourhood of \mathbf{x}^* , without further progress.
3. If gradient computation is noiseless ($\sigma = 0$), then linear convergence to optimal point.
4. A smaller stepsize t yield better converging points.

Definition 37. The **batch gradient descent** algorithm is defined as

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t^k \nabla g(\mathbf{x}^k) = \mathbf{x}^k - \frac{t^k}{|K|} \sum_{i \in K} \nabla Q_i(\mathbf{x}^k),$$

where K denotes a set of p randomly selected datapoints.

Chapter 5

Convexity

5.1 Convex Sets

Definition 38. A set $C \subseteq \mathbb{R}^n$ is called **convex** if

$$\forall \mathbf{x}, \mathbf{y} \in C \text{ and } \lambda \in [0, 1], \lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in C.$$

Equivalently, for any $\mathbf{x}, \mathbf{y} \in C$, the line segment $[\mathbf{x}, \mathbf{y}]$ is also in C .

Example 39. Very important convex sets

- A line in \mathbb{R}^n is a set of the form

$$L = \{\mathbf{z} + t\mathbf{d} : t \in \mathbb{R}\},$$

where $\mathbf{z}, \mathbf{d} \in \mathbb{R}^n$ and $\mathbf{d} \neq \mathbf{0}$.

- $[\mathbf{x}, \mathbf{y}]$, (\mathbf{x}, \mathbf{y}) for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ($\mathbf{x} \neq \mathbf{y}$), \emptyset , and \mathbb{R}^n .
- A **hyperplane** is a set of the form

$$H = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} = b\} \quad (\mathbf{a} \in \mathbb{R}^n \setminus \{\mathbf{0}\}, b \in \mathbb{R})$$

- The associated **half space** is the set

$$H^- = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} \leq b\}.$$

- The open ball $B(\mathbf{c}, r)$ and the closed ball $B[\mathbf{c}, r]$.
- The **ellipsoid** is a set of the form

$$E = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x}^T Q \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c \leq 0\}$$

where $Q \in \mathbb{R}^{n \times n}$ is positive semidefinite, $\mathbf{b} \in \mathbb{R}^n$ and $c \in \mathbb{R}$.

Lemma 40. Let $C_i \subseteq \mathbb{R}^n$ be a convex set for any $i \in I$, where I is an index set (possibly infinite), then $\bigcap_{i \in I} C_i$ is convex.

Comments: A direct consequence of the above is that convex polytopes of the form

$$P = \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \leq \mathbf{b}\},$$

are convex since they are generated as the intersection of m half-spaces $\mathbf{a}_i^T \mathbf{x} \leq b_i$.

Theorem 41. Several important algebraic properties of convex sets:

1. Let $C_1, C_2, \dots, C_k \subseteq \mathbb{R}^n$ be convex sets and let $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}$, then the set $\mu_1 C_1 + \mu_2 C_2 + \dots + \mu_k C_k$ is convex.
2. Let $C_i \subseteq \mathbb{R}^{k_i}$, $i = 1, \dots, m$ be convex sets, then the cartesian product

$$C_1 \times C_2 \times \dots \times C_m = \{(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m) : \mathbf{x}_i \in C_i, i = 1, 2, \dots, m\}$$

is convex.

3. Let $M \subseteq \mathbb{R}^n$ be a convex set and let $A \in \mathbb{R}^{m \times n}$, then the set

$$A(M) = \{A\mathbf{x} : \mathbf{x} \in M\}$$

is convex.

4. Let $D \subseteq \mathbb{R}^m$ be convex and let $A \in \mathbb{R}^{m \times n}$, then the set

$$A^{-1}(D) = \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \in D\}$$

is convex.

5.2 Convex Hull

Definition 42. Given m points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in \mathbb{R}^n$, a **convex combination** of these m points is a vector of the form

$$\lambda_1 \mathbf{x}_1 + \lambda_2 \mathbf{x}_2 + \dots + \lambda_m \mathbf{x}_m$$

where $\lambda_i \in \mathbb{R}_+$ for $i = 1, 2, \dots, m$ and satisfy $\sum_{i=1}^m \lambda_i = 1$ ($\boldsymbol{\lambda} \in \Delta_m$).

Theorem 43. Let $C \subseteq \mathbb{R}^n$ be a convex set and let $\mathbf{x}_i \in C$ for $i = 1, 2, \dots, m$. Then for any $\boldsymbol{\lambda} \in \Delta_m$, the relation

$$\sum_{i=1}^m \lambda_i \mathbf{x}_i \in C$$

holds.

Definition 44. Let $S \subseteq \mathbb{R}^n$. The **convex hull** of S , denoted by $\text{conv}(S)$, is the set comprising all the convex combinations of vectors from S :

$$\text{conv}(S) = \left\{ \sum_{i=1}^k \lambda_i \mathbf{x}_i : \mathbf{x}_i, \mathbf{x}_2, \dots, \mathbf{x}_k \in S, \boldsymbol{\lambda} \in \Delta_k \right\}$$

Comment: $\text{conv}(S)$ is the “smallest” convex set containing S .

Theorem 45. Let $S \subseteq \mathbb{R}^n$ and let $\mathbf{x} \in \text{conv}(S)$. Then

$$\exists \boldsymbol{\lambda} \in \Delta_{n+1}, \exists \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n+1} \in S \text{ s.t. } \mathbf{x} = \sum_{i=1}^{n+1} \lambda_i \mathbf{x}_i.$$

Example 46. For $n = 2$, consider the four vectors

$$\mathbf{x}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \mathbf{x}_3 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \mathbf{x}_4 = \begin{pmatrix} 2 \\ 2 \end{pmatrix},$$

and let $\mathbf{x} \in \text{conv}(\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4\})$ be given by

$$\mathbf{x} = \frac{1}{8}\mathbf{x}_1 + \frac{1}{4}\mathbf{x}_2 + \frac{1}{2}\mathbf{x}_3 + \frac{1}{8}\mathbf{x}_4 = \begin{pmatrix} \frac{13}{8} \\ \frac{11}{8} \end{pmatrix} \implies \boldsymbol{\lambda} = \begin{pmatrix} 1/8 \\ 1/4 \\ 1/2 \\ 1/8 \end{pmatrix},$$

We can find out that

$$(\mathbf{x}_2 - \mathbf{x}_1) + (\mathbf{x}_3 - \mathbf{x}_1) - (\mathbf{x}_4 - \mathbf{x}_1) = 0 \implies \boldsymbol{\mu} = \begin{pmatrix} -1 \\ 1 \\ 1 \\ -1 \end{pmatrix}.$$

Since we need to satisfy that $\forall i \in \{1, 2, 3, 4\}, \lambda_i + \alpha \mu_i \geq 0$, we need to compute

$$\epsilon = \min_{i: \mu_i < 0} \left\{ -\frac{\lambda_i}{\mu_i} \right\}$$

so that $\lambda_j + \epsilon \mu_j = 0$ for $j \in \argmin_{i: \mu_i < 0} \left\{ -\frac{\lambda_i}{\mu_i} \right\}$, thereby reducing the number of \mathbf{x}_i 's required for expressing \mathbf{x} . From the four inequalities, we can obtain that

$$\begin{cases} \alpha \leq 1/8 \\ \alpha \geq -1/4 \\ \alpha \geq -1/2 \\ \alpha \leq 1/8 \end{cases}$$

and $\epsilon = \frac{1}{8}$. Substituting $\alpha = \epsilon$, we can obtain that

$$\mathbf{x} = \frac{3}{8}\mathbf{x}_2 + \frac{5}{8}\mathbf{x}_3.$$

Definition 47. Let $S \subseteq \mathbb{R}^n$ be a convex set. A point $\mathbf{x} \in S$ is called an **extreme point** of S if $\nexists \mathbf{x}_1, \mathbf{x}_2 \in S (\mathbf{x}_1 \neq \mathbf{x}_2 \text{ and } \lambda \in (0, 1), \text{ s.t. } \mathbf{x} = \lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2$. The set of extreme point is denoted by $\text{ext}(S)$.

Theorem 48. Let $S \subseteq \mathbb{R}^n$ be a compact convex set. Then

$$S = \text{conv}(\text{ext}(S)).$$

5.3 Convex Functions

Definition 49. A function $f : C \rightarrow \mathbb{R}$ defined on a convex set $C \subseteq \mathbb{R}^n$ is called **convex** (or convex over C) if

$$\forall \mathbf{x}, \mathbf{y} \in C, \lambda \in [0, 1], f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})$$

Definition 50. A function $f : C \rightarrow \mathbb{R}$ defined on a convex set $C \subseteq \mathbb{R}^n$ is called **strict convex** if

$$\forall \mathbf{x} \neq \mathbf{y} \in C, \lambda \in (0, 1), f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) < \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})$$

Definition 51. A function is called **concave** if $-f$ is convex. Similarly, f is called **strictly concave** if $-f$ is strictly convex.

Example 52. Several examples of convex functions:

- Affine functions: $f(\mathbf{x}) = a^T \mathbf{x} + b$, where $a \in \mathbb{R}^n$ and $b \in \mathbb{R}$. Take $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\lambda \in [0, 1]$, then

$$f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) = \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}).$$

- Norms: $g(\mathbf{x}) = \|\mathbf{x}\|$. Take $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\lambda \in [0, 1]$, then

$$g(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \|\lambda\mathbf{x}\| + \|(1 - \lambda)\mathbf{y}\| = \lambda g(\mathbf{x}) + (1 - \lambda)g(\mathbf{y})$$

(Jensen's Inequality) Theorem 53. Let $f : C \rightarrow \mathbb{R}$ be a convex function where $C \subseteq \mathbb{R}^n$ is a convex set. Then $\forall \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k \in C$ and $\boldsymbol{\lambda} \in \Delta_k$,

$$f\left(\sum_{i=1}^k \lambda_i \mathbf{x}_i\right) \leq \sum_{i=1}^k \lambda_i f(\mathbf{x}_i).$$

5.4 First-order Characterization of Convex Functions

Theorem 54. Let $f : C \rightarrow \mathbb{R}$ be a continuously differentiable function defined on a convex set $C \subseteq \mathbb{R}^n$. Then

$$f \text{ is convex over } C \iff \forall \mathbf{x}, \mathbf{y} \in C, f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \leq f(\mathbf{y})$$

An analogous result holds for strictly convex functions with a strict inequality.

Comment: For a convex function f defined on \mathbb{R}^2 , the tangent plane at every point is always below f .

(Global optimality test for convex(concave) function) Theorem 55. Let f be a continuously differentiable function which is convex over a convex set $C \subseteq \mathbb{R}^n$. Then

$$\nabla f(\mathbf{x}^*) = 0 \text{ for some } \mathbf{x}^* \in C \implies \mathbf{x}^* \text{ is the global } \underline{\text{minimizer}} \text{ of } f \text{ over } C.$$

This is the same for concave function being related to global maximizer.

(Convexity of quadratic function) Theorem 56. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be the quadratic function given by $f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c$ where $A \in \mathbb{R}^{n \times n}$ is symmetric, $\mathbf{b} \in \mathbb{R}^n$, and $c \in \mathbb{R}$. Then

$$f \text{ is (strictly) convex} \iff A \succeq 0 (A \succ 0).$$

(Monotonicity of the gradient) Theorem 57. Suppose that f is a continuously differentiable function over a convex set $C \subseteq \mathbb{R}^n$, then

$$f \text{ is convex over } C \iff \forall \mathbf{x}, \mathbf{y} \in C, (\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^T(\mathbf{x} - \mathbf{y}) \geq 0.$$

An analogous result holds for strictly convex functions with a strict inequality.

Proof. If f is convex, then

$$f(y) \geq f(x) + \nabla f(x) \cdot (y - x)$$

and

$$f(x) \geq f(y) + \nabla f(y) \cdot (x - y)$$

so that by adding the above inequalities, we obtain the result. \square

5.5 Second-order Characterization of Convex Functions

Theorem 58. Let f be a twice continuously differentiable function over an open convex set $C \subseteq \mathbb{R}^n$. Then

$$f \text{ is convex over } C \iff \forall \mathbf{x} \in C, \nabla^2 f(\mathbf{x}) \succeq 0$$

Example 59. Convexity of the log-sum-exp function

$$f(\mathbf{x}) = \log(e^{x_1} + e^{x_2} + \cdots + e^{x_n}), \quad \mathbf{x} \in \mathbb{R}^n.$$

The gradient is given by

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}, \quad i = 1, 2, \dots, n.$$

Therefore, the Hessian is computed as

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}) = \begin{cases} -\frac{e^{x_i} e^{x_j}}{(\sum_{j=1}^n e^{x_j})^2} & i \neq j \\ -\frac{e^{x_i} e^{x_j}}{(\sum_{j=1}^n e^{x_j})^2} + \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} & i = j \end{cases}$$

We can thus write the Hessian matrix as

$$\nabla^2 f(\mathbf{x}) = \text{diag}(\mathbf{w}) - \mathbf{w}\mathbf{w}^T, \quad \text{with} \quad \mathbf{w} = \left(\frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \right)_{i=1}^n \in \Delta_n.$$

For any $\mathbf{v} \in \mathbb{R}^n$,

$$\mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v} = \sum_{i=1}^n w_i v_i^2 - (\mathbf{v}^T \mathbf{w})^2 \geq 0,$$

since defining $s_i = \sqrt{w_i} v_i, t_i = \sqrt{w_i}$, we have

$$(\mathbf{v}^T \mathbf{w})^2 = (\mathbf{s}^T \mathbf{t})^2 \leq \|\mathbf{s}\|^2 \|\mathbf{t}\|^2 = \left(\sum_{i=1}^n w_i v_i^2 \right) \left(\sum_{i=1}^n w_i \right) = \sum_{i=1}^n w_i v_i^2.$$

Thus $\nabla^2 f(\mathbf{x}) \succeq 0$ and hence f is convex over \mathbb{R}^n .

5.6 More Results of Convex Function

Theorem 60. Let f, f_1, f_2, \dots, f_p be convex functions over a convex set $C \subseteq \mathbb{R}^n$.

- Let $\alpha \geq 0$, then αf is a convex function over C .
- The sum function $\sum_{i=1}^p f_i$ is convex over C .
- Let $A \in \mathbb{R}^{n \times m}$ and $\mathbf{b} \in \mathbb{R}^n$. Then the function $g(\mathbf{y}) = f(A\mathbf{y} + \mathbf{b})$ is convex over the convex set $D = \{\mathbf{y} \in \mathbb{R}^m : A\mathbf{y} + \mathbf{b} \in C\}$.
- Let $g : I \rightarrow \mathbb{R}$ be a nondecreasing convex function over the interval $I \subseteq \mathbb{R}$. Assume that the image of C under f is contained in I : $f(C) \subseteq I$, then the composition of g and f defined by $h(\mathbf{x}) \equiv g(f(\mathbf{x}))$ is convex over C .

(Point-wise maximum of convex functions) Theorem 61. Let $f_1, f_2, \dots, f_p : C \rightarrow \mathbb{R}$ be p convex functions over the convex set $C \subseteq \mathbb{R}^n$, then the maximum function

$$f(\mathbf{x}) \equiv \max_{i=1,2,\dots,p} \{f_i(\mathbf{x})\}$$

is convex over C .

Theorem 62. Let $f : C \times D \rightarrow \mathbb{R}$ be a convex function defined over the set $C \times D$ where $C \subseteq \mathbb{R}^m$ and $D \subseteq \mathbb{R}^n$ are convex sets. Let

$$g(\mathbf{x}) = \min_{\mathbf{y} \in D} f(\mathbf{x}, \mathbf{y}), \quad \mathbf{x} \in C$$

where we assume that the minimum is finite. Then g is convex over C .

Example 63. The distance function from a convex set $d_C(\mathbf{x}) \equiv \inf_{\mathbf{y} \in C} \|\mathbf{x} - \mathbf{y}\|$.

Theorem 64. Let $f : C \rightarrow \mathbb{R}$ be a convex function defined over a convex set $C \subseteq \mathbb{R}^n$. Let $\mathbf{x}_0 \in \text{int}(C)$. Then $\exists \epsilon > 0, L > 0$ s.t. $B[\mathbf{x}_0, \epsilon] \subseteq C$ and

$$\forall \mathbf{x} \in B[\mathbf{x}_0, \epsilon], |f(\mathbf{x}) - f(\mathbf{x}_0)| \leq L \|\mathbf{x} - \mathbf{x}_0\|.$$

Theorem 65. Let $f : C \rightarrow \mathbb{R}$ be a convex function over the convex set $C \subseteq \mathbb{R}^n$. Let $\mathbf{x} \in \text{int}(C)$. Then

$$\forall \mathbf{d} \neq \mathbf{0}, \exists f'(\mathbf{x}; \mathbf{d}).$$

Theorem 66. Let $f : C \rightarrow \mathbb{R}$ be convex and non-constant over the nonempty convex set $C \subseteq \mathbb{R}^n$. Then f does not attain a maximum at a point in $\text{int}(C)$.

Theorem 67. Let $f : C \rightarrow \mathbb{R}$ be convex over the nonempty convex and compact set $C \subseteq \mathbb{R}^n$. Then there exists at least one maximizer of f over C that is an extreme point of C .

Chapter 6

Convex Optimization

6.1 Problem Definition

A **convex optimization problem** is a problem consisting of minimizing a convex function $f(\mathbf{x})$ over a convex set C :

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in C \end{array} \quad (\text{CVX})$$

A functional form of a convex problem can be written as

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \\ & h_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, p, \end{array}$$

where $f, g_1, g_2, \dots, g_m : \mathbb{R}^n \rightarrow \mathbb{R}$ are convex functions, and $h_1, h_2, \dots, h_p : \mathbb{R}^n \rightarrow \mathbb{R}$ are affine functions. The functional form does fit into the general formulation (CVX).

Theorem 68. Let $f : C \rightarrow \mathbb{R}$ be a (strict) convex function defined on the convex set $C \subseteq \mathbb{R}^n$. Let $\mathbf{x}^* \in C$ be a local minimum of f over C . Then \mathbf{x}^* is a strict global minimum of f over C .

Theorem 69. Let $f : C \rightarrow \mathbb{R}$ be a (strict) convex function defined on the convex set $C \subseteq \mathbb{R}^n$. Then the set of optimal solutions of the problem

$$\min \{f(\mathbf{x}) : \mathbf{x} \in C\}$$

is convex. If, in addition, f is strictly convex over C , then there exists at most one optimal solution of the problem.

Example 70.

- A convex problem:

$$\begin{array}{ll} \min & -2x_1 + x_2 \\ \text{s.t.} & x_1^2 + x_2^2 \leq 3 \end{array}$$

- A nonconvex problem:

$$\begin{aligned} \min \quad & x_1^2 - x_2 \\ \text{s.t.} \quad & x_1^2 + x_2^2 = 3 \end{aligned}$$

- **Linear Programming:**

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x} \\ (\text{LP}) : \quad & \text{s.t.} \quad A\mathbf{x} \leq \mathbf{b} \\ & B\mathbf{x} = \mathbf{g} \end{aligned}$$

- Convex quadratic problems: minimizing a convex function quadratic function subject to affine constraints. The general form is

$$\begin{aligned} \min \quad & \mathbf{x}^T Q \mathbf{x} + 2\mathbf{b}^T \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} \leq \mathbf{c} \end{aligned}$$

where $Q \in \mathbb{R}^{n \times n}$ is positive semidefinite, $\mathbf{b} \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$, $\mathbf{c} \in \mathbb{R}^m$.

6.2 Stationarity

Consider the constrained optimization problem given by

$$\min_{\mathbf{x}} \quad \{f(\mathbf{x}) : \mathbf{x} \in C\}, \quad (\text{P})$$

where $C \subseteq \mathbb{R}^n$ is closed and convex, and f is continuously differentiable over C , not necessarily convex.

Definition 71. \mathbf{x}^* is called a **stationary point** of (P) if

$$\forall \mathbf{x} \in C, \nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) \geq 0.$$

Theorem 72. Let f be a continuously differentiable function over a nonempty closed convex set C , and let \mathbf{x}^* be a local minimum of (P), then \mathbf{x}^* is a stationary point of (P).

Feasible Set	Explicit Stationarity Condition
\mathbb{R}^n	$\nabla f(\mathbf{x}^*) = \mathbf{0}$
\mathbb{R}_+^n	$\frac{\partial f}{\partial x_i}(\mathbf{x}^*) \begin{cases} = 0 & x_i^* > 0 \\ \geq 0 & x_i^* = 0 \end{cases}$
$\{\mathbf{x} \in \mathbb{R}^n : \mathbf{e}^T \mathbf{x} = 1\}$	$\frac{\partial f}{\partial x_1}(\mathbf{x}^*) = \dots = \frac{\partial f}{\partial x_n}(\mathbf{x}^*)$
$B[\mathbf{0}, 1]$	$\nabla f(\mathbf{x}^*) = \mathbf{0}$ or $\ \mathbf{x}^*\ = 1$ and $\exists \lambda \leq 0 : \nabla f(\mathbf{x}^*) = \lambda \mathbf{x}^*$

Theorem 73. Let f be a continuously differentiable function over a nonempty closed convex set $C \subseteq \mathbb{R}^n$, then \mathbf{x}^* is a stationary point of (P) iff \mathbf{x}^* is an optimal solution of (P).

6.3 Orthogonal Projection Operator

Definition 74. Given a nonempty closed convex set C , the orthogonal projection operator $P_C : \mathbb{R}^n \rightarrow C$ is defined by

$$P_C(\mathbf{x}) = \operatorname{argmin} \left\{ \|\mathbf{y} - \mathbf{x}\|^2 : \mathbf{y} \in C \right\}.$$

Theorem 75. Let $C \subseteq \mathbb{R}^n$ be a nonempty closed convex set, then

$$\forall \mathbf{x} \in \mathbb{R}^n, \exists! P_C(\mathbf{x})$$

Theorem 76. Let C be a nonempty closed convex set and let $\mathbf{x} \in \mathbb{R}^n$, then

$$\mathbf{z} = P_C(\mathbf{x}) \iff \forall \mathbf{y} \in C, (\mathbf{x} - \mathbf{z})^T(\mathbf{y} - \mathbf{z}) \leq 0.$$

Example 77.

- For $C = \mathbb{R}_+^n$,

$$P_{\mathbb{R}_+^n}(\mathbf{x}) = [\mathbf{x}]_+$$

where $[\mathbf{x}]_+ = (\max\{x_1, 0\}, \max\{x_2, 0\}, \dots, \max\{x_n, 0\})^T$.

- A **box** is a subset of \mathbb{R}^n of the form

$$B = [l_1, u_1] \times [l_2, u_2] \times \dots \times [l_n, u_n] = \{\mathbf{x} \in \mathbb{R}^n : l_i \leq x_i \leq u_i\},$$

where $l_i \leq u_i \forall i = 1, 2, \dots, n$. For this set

$$[P_B(\mathbf{x})]_i = \begin{cases} u_i & x_i \geq u_i \\ x_i & l_i < x_i < u_i \\ l_i & x_i \leq l_i. \end{cases}$$

- For the closed ball in \mathbb{R}^n , $C = B[\mathbf{0}, r]$, it holds

$$P_{B[\mathbf{0}, r]}(\mathbf{x}) = \begin{cases} \mathbf{x} & \|\mathbf{x}\| \leq r \\ r \frac{\mathbf{x}}{\|\mathbf{x}\|} & \|\mathbf{x}\| > r. \end{cases}$$

Theorem 78. Let f be a continuously differentiable function over the nonempty closed convex set C , and let $s > 0$. Then

$$\mathbf{x}^* \text{ is a stationary point of (P)} \iff \mathbf{x}^* = P_C(\mathbf{x}^* - s \nabla f(\mathbf{x}^*)).$$

6.4 Gradient Projection Method

Definition 79. We can define the **gradient mapping** as

$$G_L(\mathbf{x}) = L \left[\mathbf{x} - P_C \left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \right) \right]$$

where $L > 0$.

Comments:

- In the unconstrained case, $G_L(\mathbf{x}) = \nabla f(\mathbf{x})$ since projecting onto the same point. Otherwise, $G_L(\mathbf{x}) = \mathbf{0} \iff \mathbf{x}$ is a stationary point. We can thus consider $\|G_L(\mathbf{x})\|^2$ to be *optimality measure*.
- Slight modification can be made for gradient descent method to become **gradient projection method** for solving convex optimization problem, whose descent step becomes

$$\mathbf{x}^{k+1} = P_C \left(\mathbf{x}^k - t^k \nabla f(\mathbf{x}^k) \right)$$

to ensure that at each iteration i , x^i is within the convex set.

Theorem 80. Let $\{\mathbf{x}^k\}$ be the sequence generated by the gradient projection method for solving problem (P) with either a constant stepsize $\bar{t} \in (0, \frac{2}{L})$, where L is a Lipschitz constant of ∇f , or a backtracking stepsize strategy. Assume that f is bounded below, then:

1. The sequence $\{\mathbf{x}^k\}$ is nonincreasing.
2. $G_d(\mathbf{x}^k) \rightarrow 0$ as $k \rightarrow \infty$, where

$$d = \begin{cases} 1/\bar{t} & \text{constant stepsize} \\ 1/s & \text{backtracking.} \end{cases}$$

6.5 Separation Theorem

Definition 81. A hyperplane

$$H = \left\{ \mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} = b \right\} \quad (\mathbf{a} \in \mathbb{R}^n \setminus \{\mathbf{0}\}, b \in \mathbb{R})$$

is said to **strictly separate** a point $\mathbf{y} \notin S$ from S if

$$\mathbf{a}^T \mathbf{y} > b \quad \text{and} \quad \forall \mathbf{x} \in S, \mathbf{a}^T \mathbf{x} \leq b.$$

Theorem 82. Let $C \subseteq \mathbb{R}^n$ be a nonempty closed convex set, and let $\mathbf{y} \notin C$. Then $\exists \mathbf{p} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ and $\alpha \in \mathbb{R}$ s.t. $\mathbf{p}^T \mathbf{y} > \alpha$ and $\mathbf{p}^T \mathbf{x} \leq \alpha$ for all $\mathbf{x} \in C$.

(Farkas') Lemma 83. Let $\mathbf{c} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{m \times n}$. Then exactly one of the following systems has a solution:

- (i) $A\mathbf{x} \leq \mathbf{0}, \mathbf{c}^T \mathbf{x} > 0$.
- (ii) $A^T \mathbf{y} = \mathbf{c}, \mathbf{y} \geq 0$.

Lemma 84. Let $\mathbf{c} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{m \times n}$. Then the following two claims are equivalent:

- (a) The implication $A\mathbf{x} \leq \mathbf{0} \Rightarrow \mathbf{c}^T \mathbf{x} \leq 0$ holds true.
- (b) $\exists \mathbf{y} \in \mathbb{R}_+^m$ s.t. $A^T \mathbf{y} = \mathbf{c}$.

Theorem 85. Let $A \in \mathbb{R}^{m \times n}$, then exactly one of the following two systems has a solution:

- (i) $A\mathbf{x} < \mathbf{0}$
- (ii) $\mathbf{p} \neq \mathbf{0}, A^T \mathbf{p} = \mathbf{0}, \mathbf{p} \geq \mathbf{0}$.

6.6 KKT Conditions

Theorem 86. Consider the minimization problem

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{subject to} \quad & \mathbf{a}_i^T \mathbf{x} \leq b_i \quad i = 1, 2, \dots, m \end{aligned} \tag{LCP}$$

where f is continuously differentiable over \mathbb{R}^n , $\mathbf{a}_i \in \mathbb{R}^n$, $b_i \in \mathbb{R}$, and let \mathbf{x}^* be a local minimum point of (LCP). Then $\exists \lambda_i \geq 0$ s.t.

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \mathbf{a}_i = \mathbf{0} \quad \text{and} \quad \lambda_i (\mathbf{a}_i^T \mathbf{x}^* - b_i) = 0$$

where the above equations are called the **KKT condition** or **KKT system**.

Theorem 87. Consider the problem (LCP) where additionally f is convex. Then \mathbf{x}^* is an optimal solution $\iff \exists \lambda_i \geq 0$ s.t.

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \mathbf{a}_i = \mathbf{0} \quad \text{and} \quad \lambda_i (\mathbf{a}_i^T \mathbf{x}^* - b_i) = 0$$

Theorem 88. Consider the minimization problem

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{subject to} \quad & \mathbf{a}_i^T \mathbf{x} \leq b_i \quad i = 1, 2, \dots, m \\ & \mathbf{c}_j^T \mathbf{x} = d_j, \quad j = 1, 2, \dots, p \end{aligned} \quad (\text{LCPI})$$

where f is continuously differentiable, $\mathbf{a}_i, \mathbf{c}_j \in \mathbb{R}^n$, $b_i, d_j \in \mathbb{R}$.

- (i) (Necessity of the KKT condition) If \mathbf{x}^* is a local minimum of (LCPI), then \mathbf{x}^* satisfies the KKT condition, i.e. $\exists \lambda_i \geq 0$ and $\mu_i \in \mathbb{R}$ s.t.

$$\begin{aligned} \nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \mathbf{a}_i + \sum_{j=1}^p \mu_j \mathbf{c}_j &= 0, \\ \lambda_i (\mathbf{a}_i^T \mathbf{x}^* - b_i) &= 0, \quad i = 1, 2, \dots, m. \end{aligned}$$

- (ii) (Sufficiency in the convex case) If f is convex over \mathbb{R}^n and \mathbf{x}^* is a feasible solution of (LCPI) for which $\exists \lambda_i \geq 0$ and $\mu_i \in \mathbb{R}$ s.t. the KKT conditions are satisfied, then \mathbf{x}^* is an optimal solution of (LCPI).

Example 89. Solve the problem

$$\begin{aligned} \min \quad & \frac{1}{2}(x_1^2 + x_2^2 + x_3^2) \\ \text{s.t.} \quad & x_1 + x_2 + x_3 = 3. \end{aligned}$$

Since the function is convex, and the constraint is a hyperplane which is a convex set, KKT conditions are necessary and sufficient for this problem.

Now we assemble the **Lagrangian**

$$L(\mathbf{x}, \mu) = \frac{1}{2}(x_1^2 + x_2^2 + x_3^2) + \mu(x_1 + x_2 + x_3 - 3)$$

and solve for the KKT system, which is to find \mathbf{x}^*, μ^* s.t.

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \mu^*) = 0 \quad \text{subject to} \quad x_1 + x_2 + x_3 = 3.$$

This translates to

$$\begin{cases} \frac{\partial L}{\partial x_1} = x_1 + \mu = 0 \\ \frac{\partial L}{\partial x_2} = x_2 + \mu = 0 \\ \frac{\partial L}{\partial x_3} = x_3 + \mu = 0 \\ x_1 + x_2 + x_3 = 3 \end{cases} \implies \begin{cases} \mu = -1 \\ x_1 = x_2 = x_3 = 1. \end{cases}$$

Thus, $\mathbf{1}$ is the solution to the problem.

6.7 Orthogonal projections

Theorem 90. Let C be the affine space $C = \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}\}$ where $A \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, then

$$P_C(\mathbf{y}) = \mathbf{y} - A^T(AA^T)^{-1}(A\mathbf{y} - \mathbf{b}).$$

Example 91. Consider the hyperplane

$$H = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} = b\} \quad (\mathbf{0} \neq \mathbf{a} \in \mathbb{R}^n, b \in \mathbb{R}),$$

then by the projection on an affine space result (as shown previously),

$$P_H(\mathbf{y}) = \mathbf{y} - \mathbf{a}(\mathbf{a}^T \mathbf{a})^{-1}(\mathbf{a}^T \mathbf{y} - b) = \mathbf{y} - \frac{\mathbf{a}^T \mathbf{y} - b}{\|\mathbf{a}\|^2} \mathbf{a}.$$

Thus the distance of \mathbf{y} to the hyperplane H is

$$d(\mathbf{y}, H) = \frac{|\mathbf{a}^T \mathbf{y} - b|}{\|\mathbf{a}\|}.$$

Lemma 92. Let $H^- = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} \leq b\}$ where $\mathbf{0} \neq \mathbf{a} \in \mathbb{R}^n$ and $b \in \mathbb{R}$. Then

$$P_{H^-}(\mathbf{x}) = \mathbf{x} - \frac{[\mathbf{a}^T \mathbf{x} - b]_+}{\|\mathbf{a}\|^2} \mathbf{a}.$$

Theorem 93. Let \mathbf{x}^* be a feasible solution of

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \\ & h_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, p \end{aligned} \tag{NLP}$$

where f, g_i are continuously differentiable functions over \mathbb{R}^n and h_i are affine functions. We can define the **Lagrangian**

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^p \mu_j h_j(\mathbf{x})$$

s.t.

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}, \boldsymbol{\mu}) = 0 \quad \text{and} \quad \lambda_i g_i(\mathbf{x}^*) = 0,$$

then \mathbf{x}^* is an optimal solution of (NLP).

Theorem 94. Let \mathbf{x}^* be an optimal solution of the problem

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \\ & h_j(\mathbf{x}) \leq 0, \quad j = 1, 2, \dots, p \\ & s_k(\mathbf{x}) = 0, \quad k = 1, 2, \dots, q \end{aligned} \tag{NLP2}$$

where f, g_i are continuously differentiable convex functions over \mathbb{R}^n , h_i, s_i are affine functions. Suppose $\exists \hat{\mathbf{x}}$ satisfying the generalized **Slater's condition**:

$$\begin{aligned} g_i(\hat{\mathbf{x}}) &< 0, \quad i = 1, 2, \dots, m \\ h_j(\hat{\mathbf{x}}) &\leq 0, \quad j = 1, 2, \dots, p \\ s_k(\hat{\mathbf{x}}) &= 0, \quad k = 1, 2, \dots, q \end{aligned}$$

then $\exists \lambda_i, \eta_j \geq 0$ and $\mu_k \in \mathbb{R}$ s.t.

$$\begin{aligned} \nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x}^*) + \sum_{j=1}^p \eta_j \nabla h_j(\mathbf{x}^*) + \sum_{k=1}^q \mu_k \nabla s_k(\mathbf{x}^*) &= 0, \\ \lambda_i g_i(\mathbf{x}^*) &= 0, \quad i = 1, 2, \dots, m \\ \eta_j h_j(\mathbf{x}^*) &= 0, \quad j = 1, 2, \dots, p. \end{aligned}$$

Chapter 7

Duality

7.1 The Primal and Dual Problems

Consider the problem

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \\ & h_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, p \\ & \mathbf{x} \in X, \end{aligned} \tag{Primal}$$

where f, g_i, h_j are functions defined on the set $X \subseteq \mathbb{R}^n$. This is the “usual” optimization problem, and we refer to it as the **primal** problem. The associated Lagrangian is

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^p \mu_j h_j(\mathbf{x}) \quad (\mathbf{x} \in X, \boldsymbol{\lambda} \in \mathbb{R}_+^m, \boldsymbol{\mu} \in \mathbb{R}^p).$$

The **dual** objective function $q : \mathbb{R}_+^m \times \mathbb{R}^p \rightarrow \mathbb{R} \cup \{-\infty\}$ is defined to be

$$q(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$$

The domain of the dual objective function is

$$\text{dom}(q) = \{(\boldsymbol{\lambda}, \boldsymbol{\mu}) \in \mathbb{R}_+^m \times \mathbb{R}^p : q(\boldsymbol{\lambda}, \boldsymbol{\mu}) > -\infty\}.$$

The **dual problem** is given by

$$\begin{aligned} \max \quad & q(\boldsymbol{\lambda}, \boldsymbol{\mu}) \\ \text{s.t.} \quad & (\boldsymbol{\lambda}, \boldsymbol{\mu}) \in \text{dom}(q) \end{aligned} \tag{Dual}$$

Theorem 95. Consider the primal problem (Primal) with f, g_i, h_j and q the dual function defined in (Dual). Then

- (a) $\text{dom}(q)$ is a convex set
- (b) q is a concave function over $\text{dom}(q)$

7.2 Weak and Strong Duality

(Weak Duality) Theorem 96. Consider the primal problem (Primal) and its dual problem (Dual). Then

$$q^* \leq f^*$$

where f^*, q^* are the primal and dual optimal values respectively.

Example 97. Consider the problem

$$\begin{aligned} \min \quad & x_1^2 - 3x_2^2 \\ \text{s.t.} \quad & x_1 = x_2^3 \end{aligned}$$

We can easily solve this to obtain $\mathbf{x}^* = (1, 1)$ or $(-1, -1)$ so that $f^* = -2$.

To solve the dual problem, we formulate the Lagrangian

$$L(x_1, x_2, \lambda) = x_1^2 - 3x_2^2 + \lambda(x_1 - x_2^3)$$

so that

$$q(\lambda) = \min_{x_1, x_2} x_1^2 + \lambda x_1 - 3x_2^2 - \lambda x_2^3 = -\infty$$

since we can set $\lambda = \infty$ if $x_2 > 0$ or set $\lambda = -\infty$ if $x_2 < 0$.

Theorem 98. Consider the optimization problem f^*

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \\ & \mathbf{x} \in X \end{aligned}$$

where X is a convex set and f, g_i are convex functions over X . Suppose that $\exists \mathbf{x} \in X$ for which $g_i(\hat{\mathbf{x}}) < 0$. If this problem has a finite optimal value, then

1. the optimal value of the dual problem is obtained
2. the primal and dual problems have the same optimal value, $f^* = q^*$.

Example 99. Consider the problem

$$\begin{aligned} \min \quad & x_1^2 - x_2 \\ \text{s.t.} \quad & x_2^2 \leq 0 \end{aligned}$$

We can deduce from $x_2^2 \leq 0$ that $x_2 = 0$, and since the minimal value of x_1^2 is 0, we can deduce that $x_1 = 0$ as well, so $\mathbf{x}^* = (0, 0)$ with $f^* = 0$.

To look at the dual problem, we construct the Lagrangian

$$L(x_1, x_2, \lambda) = x_1^2 - x_2 + \lambda x_2^2$$

thereby having

$$q^* = \min_{x_1, x_2} x_1^2 - x_2 + \lambda x_2^2 = \begin{cases} -\infty & \lambda = 0 \\ -\frac{1}{4\lambda} & \lambda > 0 \end{cases}$$

and the duality problem yields 0 with $\lambda = \infty$. Note that this is however never actually attained because Slater's condition is not fulfilled.

(Complementary Slackness Conditions) Theorem 100. Consider the optimization problem

$$f^* := \min \{ f(\mathbf{x}) : g_i(\mathbf{x}) \leq 0, i = 1, 2, \dots, m, \mathbf{x} \in X \},$$

and assume that $f^* = q^*$ where q^* is the optimal value of the dual problem. Let $\mathbf{x}^*, \boldsymbol{\lambda}^*$ be feasible solutions of the primal and dual problems. Then \mathbf{x}^* and $\boldsymbol{\lambda}^*$ are optimal solutions of the primal and dual problems iff

1. $\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\lambda}^*)$
2. $\lambda_i^* g_i(\mathbf{x}^*) = 0, i = 1, 2, \dots, m.$

Comments: By establishing e.g. strong duality condition, the assumption of $f^* = g^*$ is automatically met, thereby the theorem could be applied, e.g. find $\boldsymbol{\lambda}^*$ first so that \mathbf{x}^* could be subsequently found by condition 1.

(General Strong Duality) Theorem 101. Consider the optimization problem

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \\ & h_j(\mathbf{x}) \leq 0, \quad j = 1, 2, \dots, p \\ & s_k(\mathbf{x}) = 0, \quad k = 1, 2, \dots, q \\ & \mathbf{x} \in X, \end{aligned}$$

where X is a convex set and f, g_i are convex functions over X . The functions h_j, s_k are affine functions. Suppose that $\exists \hat{\mathbf{x}} \in \operatorname{int}(X)$ for which the Slater's conditions are met. Then if the problem has a finite optimal value, then the optimal value of the dual problem

$$q^* = \max \{ q(\boldsymbol{\lambda}, \boldsymbol{\eta}, \boldsymbol{\mu}) : (\boldsymbol{\lambda}, \boldsymbol{\eta}, \boldsymbol{\mu}) \in \operatorname{dom}(q) \}$$

where

$$q(\boldsymbol{\lambda}, \boldsymbol{\eta}, \boldsymbol{\mu}) = \min_{\mathbf{x} \in X} \left[f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^p \eta_j h_j(\mathbf{x}) + \sum_{k=1}^q \mu_k s_k(\mathbf{x}) \right]$$

is attained, and $f^* = q^*$.

See the 3 examples on notes for applications of duality concept.